

Introduction

- High-quality explanations required for safety-critical ML applications
- > Despite huge popularity of rule-based and modified backpropagation eXplainable AI (XAI) methods, they often struggle with modern model architectures with innovative building blocks
- Two main problems can be identified:
 - 1) Implementation Invariance: two models implementing the same function should lead to identical explanations \rightarrow violated by many XAI methods [1]
 - 2) Parameterization: Many XAI parameter choices, optimization thereof often neglected \rightarrow non-optimal explanations

• Our contributions

-) Network Canonization for modern architectures to address implementation invariance
- 2) XAI Hyperparameter Search for optimal parameterization wrt. various XAI targets
- 3) Multi-metric **XAI Evaluation Framework** to measure improvements

Network Canonization

- **Problem**: Concatenation of linear operations hurts implementation invariance, e.g., Linear layer + BatchNorm (BN) at test time.
- > Rule-based and modified backpropagation-based approaches are affected by model structure, not only the function
- **Solution**: Re-structure model such that no neighboring linear operations exist without changing the function \longrightarrow Network Canonization [2]
- > BN parameters can be merged into neighboring Linear/Conv layers







> Straightforward for simple architectures, e.g., VGG, ResNet and EfficientNet



Complicated for highly interconnected model components, such as Dense Blocks





Optimizing Explanations by Network Canonization and Hyperparameter Search

Frederik Pahde, Galip Ümit Yolcu, Alexander Binder, Wojciech Samek, Sebastian Lapuschkin



DenseNet Canonization

- > DenseNets are highly inter**connected**: Conv layers flow into multiple BN layers and BN layers get inputs from multiple Convs
- Merging Convs/BNs is non-trivial
- **Solution:**
 - 1. Swap order of BN and modified ReLU 2. Merge BN into following Conv layer

 $\mathbf{f} \mathbf{x} \quad ext{if} (w_{ ext{BN}} > 0 ext{ and } \mathbf{x} > z)$ $ext{ThreshReLU}(\mathbf{x}) = ig \mathbf{x} \quad ext{if} \left(w_{ ext{BN}} < 0 ext{ and } \mathbf{x} < -z
ight)$ z otherwise

$$h z = \mu - rac{b_{
m BN}}{w_{
m BN}/\sqrt{\sigma+\epsilon}}$$

,
$$w_{
m BN}/\sqrt{\sigma+\epsilon}$$

Canonization Experiments

- > XAI Methods: Excitation Backprop (EB), Layer-wise Relevance Propagation (LRP)
- Dataset: ILSVRC2017 (50 random classes)
- \blacktriangleright XAI Evaluation Metrics (computed with Quantus [3]):
 - 1) Complexity: How concise are explanations?
 - 2) Faithfulness: How faithful do explanations represent the features used by the model for its prediction?
 - **3)** Localization: How well do explanations localize the object of interest?
 - 4) Robustness: How robust are explanations wrt. small input perturbations?
 - **5)** Randomization: How sensitive are explanations wrt. model randomization?

		↑ Complexity		↑ Faithfulness		↑ Local. (RRA)		↑ Local. (RMA)		↓ Robustness		\downarrow Random.	
Model	canonized	no	yes	no	yes	no	yes	no	yes	no	yes	no	yes
VGG-16	$\begin{array}{c} \text{EB} \\ \text{LRP-}\alpha 2\beta 1 \\ \text{LRP-}\varepsilon + \end{array}$	0.57 0.70 0.51	0.59 0.84 0.62	0.35 0.38 0.36	0.36 0.39 0.39	0.70 0.63 0.69	0.71 0.67 0.71	0.68 0.65 0.64	0.70 0.77 0.71	0.22 0.31 0.19	0.18 0.34 0.21	1.00 0.59 0.57	1.00 0.66 0.54
ResNet-18	$\begin{array}{l} \text{EB} \\ \text{LRP-}\alpha 2\beta 1 \\ \text{LRP-}\varepsilon + \end{array}$	0.55 0.67 0.51	0.57 0.76 0.58	0.29 0.32 0.30	0.29 0.32 0.30	0.68 0.65 0.69	0.69 0.67 0.70	0.66 0.69 0.65	0.67 0.75 0.69	0.16 0.21 0.14	0.14 0.26 0.15	0.97 0.65 0.70	0.97 0.61 0.70
EfficientNet-B0	$\begin{array}{c} \text{EB} \\ \text{LRP-}\alpha 2\beta 1 \\ \text{LRP-}\varepsilon + \end{array}$	0.85 0.75 0.50	0.70 0.77 0.73	0.24 0.29 0.28	0.27 0.20 0.30	0.73 0.72 0.75	0.67 0.65 0.75	0.79 0.79 0.69	0.72 0.73 0.79	0.42 0.48 0.12	0.33 0.49 0.21	0.99 0.57 0.61	1.00 0.51 0.65
DenseNet-121	EB LRP- $\alpha 2\beta 1$ LRP- ε +	0.66 0.82 0.67	0.62 0.81 0.66	0.15 0.25 0.26	0.31 0.33 0.33	0.58 0.64 0.70	0.72 0.71 0.74	0.53 0.68 0.71	0.73 0.81 0.77	0.57 0.65 0.63	0.17 0.28 0.19	0.75 0.40 0.39	0.89 0.44 0.48

Results for ILVCR2017 for metrics Sparseness (Complexity), Region Perturbation (Faith-fulness), Relevance Rank Accuracy and Relevance Mass Accuracy (Localization). Avg. Sensitivity (Robustness) and Random Logit Test (Randomization). Arrows indicate whether high (\uparrow) or low (\downarrow) are better. Best results are shown in bold.

Results for additional datasets (Pascal VOC, MS COCO) and task (VQA on CLEVR XAI) in paper





 \blacktriangleright Alternative: Create multiple copies of linear/conv layers (one per BN layer) \rightarrow memory-consuming



XAI Hyperparameter Search

Experiment:

- ImageNet
- LRP γ -rule as XAI method
- Parameters:

- 2. canonization yes/no



Conclusions

- Future work:

References

1] Montavon, Grégoire. "Gradient-based vs. propagation-based explanations: An axiomatic comparison." Explainable ai: Interpreting, explaining and visualizing deep learning (2019): 253-265.

[2] Motzkus, Franz, et al. "Measurably stronger explanation reliability via model canonization." IEEE International Conference on Image Processing (2022).

3] Hedström, Anna, et al. "Quantus: An Explainable AI Toolkit for Responsible Evaluation of Neural Network Explanations and Beyond." Journal of Machine Learning Research 24.34 (2023): 1-11.







> Problem: Flexibility of XAI methods, e.g., LRP, leads to numerous possible parameterizations (e.g., LRP-rule/layer-assignment, rule parameters, canonization)

Multi-metric Hyperparameter Search: Measure impact of different parameterizations on different metrics, estimating the explanation quality from various viewpoints

- VGG-13 model (with BN), pretrained on

$$R_j = \sum_k rac{a_j \cdot (w_{jk} + \gamma \, w_{jk}^+)}{\sum_j a_j \cdot (w_{jk} + \gamma \, w_{jk}^+)} \cdot R_k$$

1. γ per layer (for simplicity, we group layers into Conv 1-3, 4-7, 8-10, Classifier) with $\gamma \in \{0, 0.1, 0.25, 0.5, 1, 10\}$

 $\gamma = 0$: negative/positive contributions are treated equally

 $\gamma \rightarrow \infty$: negative contributions are neglected

Introduced model canonization for DenseNet and EfficientNet

Impact of canonization depends on the model architecture and XAI evaluation metric, but canonization provides an extra option to tune explanations for the task at hand

> XAI hyperparameter choices are key. They can be tuned using our multi-metric evaluation framework. Optimal hyperparameter choice depends on the given task.

- extend canonization to other relevant model architectures, including Vision Transformers - optimze hyperparameter search beyond grid search with multi-target objective